

Game of Duels: Information-Theoretic Axiomatization of Scoring Rules

Jakša Cvitanić¹, Dražen Prelec², Sonja Radas³, and Hrvoje Šikić⁴

Abstract — This paper aims to develop insights into Bayesian truth serum (BTS) mechanism by postulating a sequence of seven natural conditions reminiscent of axioms in information theory. The condition that reduces a larger family of mechanisms to BTS is additivity, akin to the axiomatic development of entropy. The seven conditions identify BTS as the unique scoring rule for ranking respondents in situations in which respondents are asked to choose an alternative from a finite set and provide predictions of their peers’ propensities to choose, for finite or infinite sets of respondents.

Index Terms— Bayesian Truth Serum, information entropy, Shannon theory

I. INTRODUCTION

The Bayesian truth serum (BTS) algorithm [1] is a game-theoretic scoring system, designed to incentivize honest responses to non-verifiable questions. For each multiple-choice question in a survey, the respondent is asked to both answer the question and also to predict the distribution of answers by the rest of the survey sample. The prediction is expressed in terms of percentages of respondents that will choose each possible answer. Once these two inputs are collected from all respondents, the algorithm assigns to each respondent a numerical BTS score, calculated via a mathematical formula (that we recall below).

This paper was first submitted for review on February 25, 2016.

¹Division of the Humanities and Social Sciences, Caltech. E-mail: cvitanic@hss.caltech.edu. Research supported in part by NSF grant DMS 10-08219.

²MIT, Sloan School of Management, Department of Economics, Department of Brain and Cognitive Sciences. E-mail: dprelec@mit.edu.

³The Institute of Economics, Zagreb. and MIT, Sloan School of Management. E-mail: sradas@mit.edu This research was supported by a Marie Curie International Outgoing Fellowship within the 7th European Community Framework Programme, PIOF-GA-2013-622868 - BayInno.

⁴University of Zagreb, Faculty of Science, Department of Mathematics. E-mail: hskic@math.hr. Research supported in part by the MZOS grant 037-0372790-2799 of the Republic of Croatia and in part by Croatian Science Foundation under the project 3526

This paper was presented [in part] at Bayesian Crowd Workshop, July 3-4, 2017, Erasmus University Rotterdam.

The original paper on BTS [1] defined conditions under which the scoring rule is strictly incentive-compatible, which means that an honest answer on each question strictly maximizes that respondent’s expected score, assuming that other respondents are answering honestly and the sample size can be made arbitrarily large. BTS incentives have been applied to a range of survey settings, including knowledge design [2], criminology [3], economics and psychology [4], and new product adoption [5]¹.

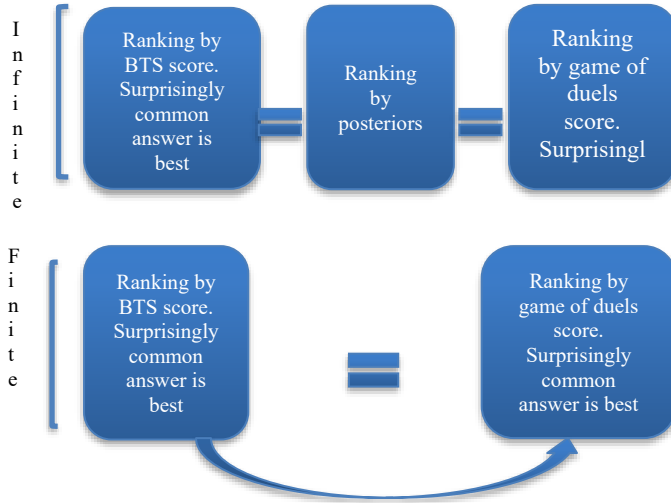
This paper is concerned with a different property of the BTS score, namely, that it generates a ranking of respondents that reflects the quality of their information, or domain expertise. We show that a finite version of the BTS score can be obtained as the outcome of a “game of duels” in which each player engages in a duel with every other player (including himself). That is, for each player, under natural conditions on the rules of the game, the payoffs in the “game of duels” are exactly those of BTS. The key condition is the additivity property as employed in the Shannon information theory.

It is known that the ranking by the BTS score, in the case of infinitely many players, corresponds to the ranking by posterior probabilities of the true state of nature, called “posteriors” (see [1] and [6]). Unfortunately, this property fails in the case when there are finitely many players. While there exist mechanisms that are incentive compatible in the finite case (e.g., [15], [16], [17] etc.), it is not difficult to show that no finite case algorithm will rank players by posteriors. In this paper we, instead, rank the players by having them compete pairwise in scored duels. This can be done both in the finite and in the infinite case. The main contribution of the paper is to identify natural conditions under which such a game reproduces BTS scores.

¹ For numerous references for the study of various truth-inducing scoring rules in the game-theoretic context with many players see [6]. When only one respondent is asked to reveal an opinion on a probability distribution, the mechanisms that incentivize truth-telling are called proper scoring rules. The literature goes back all the way back to [7], [8] and [9]. Papers that make a connection between proper scoring rules and entropy include [10], [11], [12], [13] and [14].

Let us elaborate on the connection between BTS, ranking by game of duels and ranking by posteriors, to which we refer as PstRn or posterior ranking. Recently, it was shown in [18] that with infinitely many respondents, the best PstRn expert is also the respondent who selects the answer that is most ‘surprisingly common,’ that is, most underestimated relative to predictions. Although the best expert according to PstRn cannot be identified in the finite case, we show that it is possible instead to identify the person who selects the answer that is most surprisingly common through a series of pairwise comparisons (or ‘duels’). The ranking of respondents in this contest serves as a proxy for the PstRn ranking in the finite case. Figure 1 displays the relationships.

Figure 1. Comparison of BTS ranking in finite and infinite samples



In our model, players play a series of (conceptual) duels. After each duel, points are transferred from one player to the other². A player’s final score is the total number of points received (or lost). The respondents are ranked according to their scores. The nature of the game makes it especially suitable for situations when players are machines. Although this approach seems to have little in common with BTS, our main contribution lies in establishing a connection between the two. Notably, the games of duels in which transfers satisfy certain conditions rank the players according to how “generously” they predict the shares of the answers they have not chosen, and with the additional additivity condition, the only possible game of duels is the one that results in BTS scores.

² This means that each of the duels results in a transfer of points between players. The order of duels is not important, and there is no interdependence of duels. Although for each of the players her/his duels occur in a sequence, the procedure can be implemented so that different set of players engage in duels simultaneously.

II. BAYESIAN TRUTH SERUM ALGORITHM

Here we give a short theoretical exposition of the Bayesian Truth Serum. We denote by R the set of players (respondents). We assume that R is not empty, not a singleton, and at most countable (i.e. the cardinal number of the set R satisfies $2 \leq \text{card}(R) \leq \aleph_0$). Suppose that the players are presented with a multiple choice question, offering a choice of $m \in \mathbb{N} \setminus \{1\}$ answers (we use the standard mathematical notation where \mathbb{N} is the set of natural numbers, \mathbb{R} is the set of real numbers, and $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty\} \cup \{+\infty\}$). Each player picks a simple answer (the one s/he thinks is the correct one) and gives a prediction in terms of probabilities about the distribution of m answers within R .³ More precisely, we present the answer of a player $r \in R$ as a pair of ordered probability vectors

$$((x_1^r, \dots, x_m^r); (y_1^r, \dots, y_m^r)) ; \quad (1)$$

where $x_1^r, \dots, x_m^r \in \{0,1\}$, and $y_1^r, \dots, y_m^r \in [0,1]$ such that $\sum_{k=1}^m x_k^r = 1$ and $\sum_{k=1}^m y_k^r = 1$. Exactly one of x_k^r is equal to one (the non-zero term which corresponds to the selected answer), while (y_1^r, \dots, y_m^r) is a probability distribution on $\{1, 2, \dots, m\}$. As a consequence, answers of all the players can be presented as a (finite or infinite) matrix $(X; Y)$; it is of the order $\text{card}(R) \times 2m$ and its r^{th} row, $r \in R$, is given by (1).

We want to assign a numerical score to each player based on $(X; Y)$, denoted

$$u^r = u^r(X; Y); \quad (2)$$

for player $r \in R$. Eventually, we expect our scores to be real-valued, but here at the outset we shall not restrict ourselves and in principle we allow even for infinite values, i.e.

$$u^r(X; Y) \in \overline{\mathbb{R}}. \quad (3)$$

A. The score in Bayesian Truth Serum

To develop the formula for the score, we shall use the notation $\sum_{s \in R}$ in both finite and infinite case. If R is finite, then $\sum_{s \in R}$ has its usual meaning of the sum over all elements of R . We define $\bar{x} := (\bar{x}_1, \dots, \bar{x}_m)$ where $\bar{x}_k := \frac{1}{\text{card}(R)} \sum_{s \in R} x_k^s$, for $k=1, \dots, m$. It is easy to see that \bar{x}_k represent arithmetic means of X -columns. We also define $\hat{y} := (\hat{y}_1, \dots, \hat{y}_m)$ where $\ln(\hat{y}_k) := \frac{1}{\text{card}(R)} \sum_{s \in R} \ln(y_k^s)$ for $k=1, \dots, m$. Here \hat{y}_k are geometric means of Y -columns.

If R is infinite, then we write $R = \bigcup_{n \in \mathbb{N}} R_n$, where $\text{card}(R_n) = n$, and the meaning of $\sum_{s \in R}$ is in the sense of

³ The latter question is usually asked in the following way: “please estimate the percentage of your peers who will choose answer k ”, and the question is repeated for each $k=1, \dots, m$.

$\lim_{n \rightarrow \infty} \sum_{s \in R_n}$; the notation comes together with the assumption that the limit exists within $\overline{\mathbb{R}}$. We extend the definition of $\bar{x} := (\bar{x}_1, \dots, \bar{x}_m)$ so that we define

$\bar{x}_k := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{s \in R_n} x_k^s$. Similarly, we extend the definition of $\hat{y} := (\hat{y}_1, \dots, \hat{y}_m)$ by defining $\ln(\hat{y}_k) := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{s \in R_n} \ln(y_k^s)$.

Using the notation above, the respondent's BTS score in [1] is defined as

$$u^r(X; Y) := \sum_{k=1}^m x_k^r \ln \frac{\bar{x}_k}{\hat{y}_k} + \sum_{k=1}^m \bar{x}_k \ln \frac{y_k^r}{\bar{x}_k}, \quad (4)$$

where $r \in R$. The first part of the sum is called the information score, while the second one is called the prediction score [1].

B. Bayesian Truth Serum in applications

BTS method can be applied in survey settings, as explained in [1]. In applications, this algorithm works as follows:

1. It is explained to the respondents that they will be rewarded according to the BTS scoring rule. The rule itself is not explained, except that the respondents are told that it is incentive compatible.
2. Respondents are asked to report their answer from m offered alternatives. For a chosen respondent r this will create the vector (x_1^r, \dots, x_m^r) .
3. Respondents are asked to predict how others will choose. This will create the vector (y_1^r, \dots, y_m^r) for the respondent r .
4. Respondents are rewarded according to the BTS scoring rule (outlined in (4)).

It was shown in [1] that the BTS scoring rule is budget balanced, and allows a strict Nash equilibrium in which everyone responds honestly. It is shown in [1] and [6] that rank-ordering respondents by their BTS score is the same as rank-ordering them by their posterior probability for the realized state of nature. In [19] it was experimentally demonstrated that BTS alters respondents' behavior in the desired direction, which makes it suitable for survey applications.

C. A Bayesian framework and ranking of players

For theoretical studies of BTS, the following Bayesian framework is assumed. We assume that respondents are presented via a family $(V(r) : r \in R)$ of random variables

taking values in $\{1, \dots, m\}$. We also assume that there is a random variable Ω , the actual state of nature, with finitely many values N . It is standard to assume that $(V(r) : r \in R)$ are Ω -conditionally i.i.d. Hence, the complete probability distribution of the system is given via the distribution of the 2-dimensional random vector $(V(r_0), \Omega)$ (notice that we can take any $r_0 \in R$ here due to Ω -conditional i.i.d. assumption). Obviously, this distribution is given as a $m \times N$ probability matrix Q . In particular, the probabilities

$$p_{jk} := P(V(r) = j, V(s) = k), \quad r \neq s$$

do not depend on the choice of r and s , as long as $r \neq s$; $r, s \in R$. Within this Bayesian framework the theoretical analysis of the system is done under the assumption that the values y_j^r are given through Bayesian updating (see, for example [6] for details). More precisely, assuming that $x_k^r = 1$, we have

$$y_j^r = P(V(r) = j | V(s) = k) = \frac{p_{jk}}{\sum_{l=1}^m p_{lk}}.$$

We can write $x_k^r = 1$ as $V(r) = k$. Notice that the vectors (x_k^r) and (y_j^r) allow us to compute the BTS payoff $u^r(X; Y)$ to player r . Then, it can be shown that (see [1] and [19]),

$$u^r = u^r(X; Y) = \ln(P(\Omega = i_0 | V(r) = k)) - \sum_{j=1}^m P(V(r) = j | \Omega = i_0) \ln(P(\Omega = i_0 | V(r) = j));$$

where i_0 denoted the true state of nature and $\text{card}(R) = \aleph_0$. In particular, the above formula shows that for the infinite set of respondents BTS is increasing in the posterior probabilities, a property called Posterior Ranking or PstRn. However, PstRn does not hold with finitely many players. We now identify a property that is equivalent to PstRn with infinitely many players, which will hold also in the finite case under our mechanism.

For simplicity, let us turn to the binary case where $m = 2$. We denote possible answers as Y and N . We also assume that there are two states of nature "True" and "False". We identify states of nature with distributions on (Y, N) , i.e. "True" = $(T, 1-T)$ and "False" = $(F, 1-F)$, where $T, F \in (0, 1)$ and $T \neq F$. Observe that $T = P(V(r) = Y | \Omega = \text{True})$ and analogous formula holds for F . We also denote $P(\Omega = \text{True})$ as $P(T)$ and $P(\Omega = \text{False})$ as $P(F)$.

We introduce the property of "being modest to oneself", called Mds, defined as, for player r ,

$$\frac{T}{y_Y^r} > \frac{1-T}{y_N^s};$$

where $x_Y^r = 1$, $r \neq s$ and $x_N^s = 1$. This condition essentially considers how players predict the share of their chosen

answers compared to the realized percentage. Player r satisfies Mds if she underestimates the share of her chosen answer Y more than the player s does for his chosen answer N . In that sense the player r is more “modest”.

We also introduce the property “being generous to others”, called Gen, by

$$\frac{T}{y_Y^s} > \frac{1-T}{y_N^r};$$

where $x_Y^r = 1$, $r \neq s$ and $x_N^s = 1$. While Mds considers how players predict the percentage of people who choose the same answer as them, the condition Gen considers how players predict the share of the opposite answer. Player r satisfies Gen if her prediction of the opposite answer N is closer to the real percentage compared to what player s predicts for the answer Y . In other words, the player r underestimates her non-chosen answer less than player s underestimates his non-chosen answer. In this sense player r is more “generous” than player s .

Observe that Mds and Gen can both be interpreted as “the player has selected a surprisingly common answer”.

An elementary calculation shows that we have

$$\begin{aligned} p_{YY} &= T^2 P(T) + F^2 P(F) \\ p_{YN} &= p_{NY} = T(1-T)P(T) + F(1-F)P(F) \\ p_{NN} &= (1-T)^2 P(T) + (1-F)^2 P(F). \end{aligned}$$

Furthermore, we have, for $x_Y^r = 1$,

$$y_Y^r = \frac{T^2 P(T) + F^2 P(F)}{TP(T) + FP(F)};$$

and, similarly, for $x_N^r = 1$,

$$y_N^r = \frac{(1-T)^2 P(T) + (1-F)^2 P(F)}{(1-T)P(T) + (1-F)P(F)}.$$

It is now a straightforward algebraic calculation to check that under the assumptions that $r, s \in R$, $r \neq s$, $x_Y^r = 1$, $x_N^s = 1$ we have

$$u^r > u^s \quad (5)$$

$$\Leftrightarrow P(\Omega = \text{True} | V(r) = Y) > P(\Omega = \text{True} | V(s) = N) \quad (6)$$

$$\Leftrightarrow T > F \quad (7)$$

$$\Leftrightarrow \frac{T}{y_Y^r} > \frac{1-T}{y_N^s} \Leftrightarrow \frac{T}{y_Y^s} > \frac{1-T}{y_N^r}. \quad (8)$$

Observe that the first and the second equivalence do not make sense if we are not within a stochastic framework. The first equivalence is PstRn. Hence, we have five equivalent conditions, one of which is BTS. However, in the finite case, the above equivalences do not all hold. We will show below that in our deterministic mechanism (“game of duels”) the Gen

equivalence remains and it is valid both in the finite and the infinite case.

III. A SYSTEM OF CONDITIONS

We will develop a system of conditions that results in scores (4). In our approach players get ranked via simultaneous conceptual duels. Each duel has a “challenger”, player $r \in R$, and an “offender”, player $s \in R$.⁴ We denote such duel as $r \rightarrow s$. Each respondent plays a duel with every other respondent, including oneself.

Each duel $r \rightarrow s$ ends with a transfer of points from player r to player s . We denote the number of transferred points by

$$T^{r \rightarrow s} = T^{r \rightarrow s}(X; Y) \in \mathbb{R}. \quad (9)$$

We can think of positive $T^{r \rightarrow s}$ as the winning case for the offender, while negative $T^{r \rightarrow s}$ means that the challenger prevails. All the possible duels are to be performed (including the duel with oneself) in order to determine scores u^r for all respondents $r \in R$. In particular, if R is finite, there will be $[card(R)]^2$ duels.

Let us introduce the basic rule for a duel. For every $r \in R$ the score u^r equals the number of received points minus the number of given points, i.e.

$$u^r = u^r(X, Y) = \sum_{s \in R} T^{s \rightarrow r}(X; Y) - \sum_{s \in R} T^{r \rightarrow s}(X; Y) \quad (10)$$

There are two immediate important consequences of (10). First, assuming that all the sums are finite-valued (which is the only interesting case), the duel is a zero-sum game,

$$\sum_{r \in R} u^r = \sum_{r \in R} \sum_{s \in R} T^{s \rightarrow r} - \sum_{r \in R} \sum_{s \in R} T^{r \rightarrow s} = 0. \quad (11)$$

The second consequence of (10) is that the description of u^r reduces to the description of $T^{r \rightarrow s}$. We will present a set of seven conditions about $T^{r \rightarrow s}$ that generate BTS algorithm (4). For each condition we give an intuitive justification (which may include some ideas from statistics) and a formal statement (which is always going to be deterministic).

The first six conditions are natural to impose and their combined effect will be that, for every $r, s \in R$, for some function P we have

$$T^{r \rightarrow s}(X; Y) = \sum_{k=1}^m x_k^s P(\bar{x}_k; y_k^r);$$

where \bar{x}_k is the sample mean. The seventh condition will be

⁴ We use traditional duel terminology, where one player (offender) offends the other (challenger), who in turn challenges the first player to a duel

the additivity condition, which will reduce the above representation to BTS.

Our first condition is very much in the spirit of medieval duels. We can interpret it as “the offender chooses the playground for the duel”.

Condition 1. *The challenger r will transfer points to the offender s based on the x answer of the offender s . More precisely, for every $r, s \in R$ and for every $k \in \{1, \dots, m\}$ there exists a number $P_k^{rs}(X; Y) \in \mathbb{R}$ such that*

$$T^{r \rightarrow s}(X; Y) = \sum_{k=1}^m x_k^s P_k^{rs}(X; Y). \quad (12)$$

Observe that Condition 1 reduces our analysis from $T^{r \rightarrow s}$ to P_k^{rs} . Observe also that, for every $s \in R$, there is exactly one $k \in \{1, \dots, m\}$ such that $x_k^s = 1$. Hence, we can think of that k as being the function of s , i.e. $k = k(s)$. It follows then that (12) becomes

$$T^{r \rightarrow s}(X; Y) = P_{k(s)}^{rs}(X; Y). \quad (13)$$

In order to understand the second condition, we introduce the following partition of R

$$R_k := \{s \in R \mid x_k^s = 1\}, \quad k = 1, \dots, m. \quad (14)$$

Obviously, the partition $R = R_1 \cup \dots \cup R_m$ is a function of X . Fix k for a moment and consider R_k , which is a subset of players who choose the same answer k . In general, the number of points P_k^{rs} may vary as s changes within R_k . The purpose of our second condition is to prevent this from happening, i.e. that condition can be thought of as “the egalitarian principle within R_k .”

Condition 2. *Given $r \in R$ and $k \in \{1, \dots, m\}$ we have*

$$s, s' \in R_k \Rightarrow P_k^{rs}(X; Y) = P_k^{rs'}(X; Y).$$

Condition 2 says that if offenders $s, s' \in R$ choose the same answer, then in the duels with all challengers they will receive the same number of points. Observe that Condition 2 includes even the cases when for some k the set R_k may be an empty set; in this case the implication in Condition 2 is true, since the premise of the implication is never true. Using a slight abuse of notation (think of $k=k(s)$), Condition 2 implies that

$$P_k^{rs}(X; Y) = P_k^r(X; Y). \quad (15)$$

In order to understand the third condition, observe that by choosing the answer k , the offender s decides (given that r is known) on a type of function P_k^r that will be used in the duel $r \rightarrow s$. However, the P_k^r will in general still depend on $(X; Y)$. Our next condition can be thought of as strengthening Condition 1. The offender s chooses the playground k , and in

doing so it reduces the variable dependence accordingly.

Condition 3. *For every $r \in R$ and for every $k \in \{1, \dots, m\}$,*

$$P_k^r(X; Y) = P_k^r((x_k^q)_{q \in R}; (y_k^q)_{q \in R}).$$

Next we turn to Condition 4 which has a deterministic form, but which can be justified using some ideas from statistics. One of the main problems in statistical analysis is to make inference about some unknown parameter θ . The inference is based on the information given in a sample X_1, \dots, X_n . If t is a sufficient statistic for θ , then whenever we have two sample points $x = (x_1, \dots, x_n)$ and $x' = (x'_1, \dots, x'_n)$ with the property $T(x) = T(x')$, then the inference about θ is the same regardless whether x or x' is observed. A typical example is a Bernoulli sample in which the sufficient statistics for the probability of success is the sample mean.

We argue here that the X -part of our data is akin to the Bernoulli sample set-up. We are interested in $\omega = (\omega_1, \dots, \omega_n)$, where ω_k gives the actual fraction of the population that thinks k is the correct answer to the original question. Hence, since we are interested in ω_k , then the average value gives as much information about ω_k as the entire k -th column of the matrix X , i.e. $(x_k^q)_{q \in R}$. Therefore, we term our fourth condition “the data reduction principle for X ”.

Condition 4. *For every $r \in R$ and for every $k \in \{1, \dots, m\}$,*

$$P_k^r((x_k^q)_{q \in R}; (y_k^q)_{q \in R}) = P_k^r(\bar{x}_k; (y_k^q)_{q \in R}).$$

Our second data reduction principle deals with Y . Our conditions so far provided the offender s with the advantage to “choose the playground” k . In the next condition we give an advantage to the challenger r by giving him/her an option to “choose the weapon”. We can think of it as allowing the challenger to select some information from the k^{th} column of Y in order to predict ω_k . We assume that the challenger is very self-confident and uses only his/her own choice y_k^r . This gives us the data reduction principle for Y .

Condition 5. *For every $r \in R$ and for every $k \in \{1, \dots, m\}$,*

$$P_k^r(\bar{x}_k; (y_k^q)_{q \in R}) = P_k^r(\bar{x}_k; y_k^r).$$

Observe that our conditions have reduced a function defined on a matrix $(X; Y)$ to a function defined on a pair of numbers $(\bar{x}_k; y_k^r)$ which are between 0 and 1. However, at this level of generality we still allow the form of the function to change with r or with k (i.e. the function can vary with the choice of different players or responses). A system that would allow for such level of generality would not be very practical, as for

every k and every r we would have a different function P_k^r . Hence we opt for a more robust selection and introduce the following “universality condition”.

Condition 6. *There exists a function $P: [0,1] \times [0,1] \rightarrow \mathbb{R}$ such that for every $r \in R$ and for every $k \in \{1, \dots, m\}$ we have $P_k^r = P$.*

In other words, Condition 6 ensures that function P_k^r is the same for every player r and for every answer k .

To recap, the first six conditions imply that, for every $r, s \in R$

$$T^{r \rightarrow s}(X; Y) = \sum_{k=1}^m x_k^s P(\bar{x}_k; y_k^r). \quad (16).$$

Remark on ranking of players: Consider a finite set R and a function P given by

$$P(x, y) = \frac{1}{\text{card}(R)} [f(x) - f(y)];$$

where $f: (0,1) \rightarrow \mathbb{R}$. For the purpose of this discussion, let us also assume that the same x response implies the same y response, i.e., $(x_k^r = 1 = x_k^s \Rightarrow y_j^r = y_j^s)$ for every $j \in \{1, \dots, m\}$. Then, we can use notation $y_j^k = y_j^r$ for $x_k^r = 1$. It is not difficult to calculate u^r for $x_k^r = 1$. We obtain

$$u^r = f(\bar{x}_k) - \sum_{l=1}^m \bar{x}_l f(y_l^k) - \sum_{l=1}^m \bar{x}_l (f(\bar{x}_l) - f(y_l^k)).$$

Consider now the case $m=2$, with two answers being Y and N . To simplify notation, denote \bar{x}_Y by p , y_Y^Y by y , and y_Y^N by z . If $x_Y^r = 1$, then we denote u^r by u^Y (and similarly for u^N). Observe that we are in a deterministic situation, so we do not have neither states of nature nor y_j^r which are given by Bayesian update. Hence, $y, z \in (0,1)$ with $y \neq z$ as the only requirement. We then obtain

$$\begin{aligned} u^Y &= f(p) - (pf(y) + (1-p)f(z)) - [p(f(p) - f(y)) + \\ &\quad + (1-p)(f(1-p) - f(1-y))] = \\ &= (1-p)[f(p) - f(z) + f(1-y) - f(1-p)]; \\ u^N &= f(1-p) - (pf(1-y) + (1-p)f(1-z)) - \\ &\quad - [p(f(p) - f(z)) + (1-p)(f(1-p) - f(1-z))] = \\ &= p[f(1-p) - f(1-y) + f(z) - f(p)]. \end{aligned}$$

It follows then that

$$u^Y > u^N \Leftrightarrow f(1-y) - f(1-p) > f(z) - f(p).$$

It is easier to follow the argument if we assume that f is also a strictly increasing function. Observe that the above condition is then essentially Gen-type condition, in the sense that Y player has a higher score if and only if she predicts the opposite answer more generously (in the sense of an f

increment) than N player predicts the opposite answer.

If we want to have exactly the Gen condition, then we need the “same f increments”, i.e., we need $f(x_1) - f(x_2) = f(\frac{x_1}{x_2})$. In other words, we need the additivity property. Interestingly enough, this property works even more generally, and our last condition takes this point into consideration.

Before turning back to our condition system, let us observe that in a deterministic framework, i.e., when $y, z \in (0,1)$ with $y \neq z$, conditions Gen and Mds are not equivalent. Given $p \in (0,1)$, Mds says that $p/y > (1-p)/z$, which is equivalent to $z > ((1-p)/p)y$. On the other hand, Gen says that $p/z > (1-p)/(1-y)$, which is equivalent to $z < (p/(1-p))(1-y)$.

Let us now turn our attention to the last and the most demanding condition. In order to justify it, we borrow ideas from information theory⁵. Consider two identical games of duels with the same players participating, with transfers $P(\bar{x}_k^i; y_k^{ri})$, $i=1,2$. Assume each player chooses an alternative in the second game independently from his choice in the first game, and independently of each other. Also consider a hypothetical “combined” third game that considers the pair alternatives the players have made in the first two games. Denote by \bar{x}_{kl} the proportion of the players choosing alternative (k,l) . If the number of players is large, under independence assumption we have approximately $\bar{x}_{kl} = \bar{x}_k^l \cdot \bar{x}_k^l$. Then, the additivity condition translates into a “scaling of transfers” condition: the corresponding transfers in the combined game should be equal to the sum of transfers in the two original games. In other words, if a game is composed of (independent) subgames, the transfers should scale at the same rate as the number of subgames.

As in [20] we exclude the case of zero and treat it separately (see also [1]). Hence, we introduce the additivity property condition in the following form.

Condition 7. *The restriction $P|_{(0,1] \times (0,1]}$ of the function P given in (12) is a continuous function such that, for every $u \in (0,1]$, $P(u; u) = 0$, and for every $u_1, u_2, v_1, v_2 \in (0,1]$,*

$$P(u_1 u_2; v_1 v_2) = P(u_1; v_1) + P(u_2; v_2).$$

Observe that if the selected “playground information” of the offender results in \bar{x}_k which is exactly equal to the “challenger

⁵ In particular, one may consult a chapter on a measure of information in [20] with the emphasis on section 1.2.

information”, then the natural outcome is “a draw”, i.e. $P(u, u) = 0$. As in Shannon theory, the consequence of Condition 7 is the following well known result:

Lemma. *If $h: \langle 0, 1 \rangle \rightarrow \mathbb{R}$ is continuous and such that, for every $u, v \in \langle 0, 1 \rangle$, $h(uv) = h(u) + h(v)$, then $h(u) = a \cdot \ln(u)$, where $a = -h(e^{-1})$.*

Recall that the additivity property is very strong. The conclusion of the lemma follows even with much milder requirements than continuity on function h ; for example it is sufficient to require monotonicity or measurability. Although this would allow us to reduce the requirement on continuity given in Condition 7, in order to avoid unnecessary mathematical intricacies we presented Condition 7 in the above form.

The lemma implies:

Corollary. *If a function $P: \langle 0, 1 \rangle \times \langle 0, 1 \rangle \rightarrow \mathbb{R}$ satisfies Condition 7, then there exists $a \in \mathbb{R}$ such that, for every $u, v \in \langle 0, 1 \rangle$,*

$$P(u; v) = a \cdot \ln\left(\frac{u}{v}\right).$$

Proof. Take $u_1 = u$, $u_2 = 1$, $v_1 = v$, $v_2 = 1$ in Condition 7. We obtain $P(u; v) = P(u; 1) + P(1; v)$. We start with the function $u \rightarrow P(u; 1)$. If we apply Condition 7 with $v_1 = v_2 = 1$, we obtain

$$P(u_1 u_2; 1) = P(u_1; 1) + P(u_2; 1).$$

Hence, $u \rightarrow P(u; 1)$ satisfies the requirement of the lemma. We conclude that there exists $a \in \mathbb{R}$ such that $P(u; 1) = a \cdot \ln(u)$.

Consider now the function $v \rightarrow P(1; v)$. If we apply Condition 7 with $u_1 = u_2 = 1$, we obtain

$$P(1; v_1 v_2) = P(1; v_1) + P(1; v_2).$$

Again, using the lemma we conclude that there exists $b \in \mathbb{R}$ such that $P(1; v) = b \cdot \ln(v)$.

Finally, using $P(u; u) = 0$ and $P(u; u) = P(u; 1) + P(1; u) = a \cdot \ln(u) + b \cdot \ln(u)$, we obtain $b = -a$. Hence, for every $u, v \in \langle 0, 1 \rangle$, it follows $P(u; v) = a \cdot \ln\left(\frac{u}{v}\right)$.

Q.E.D.

Remark. We need to decide on a particular choice of the normalizing constant $a \in \mathbb{R}$ from the previous corollary. Suppose for the moment that the challenger r has selected

$y_k^r = 1$, for some k . This implies $y_l^r = 0$ for all $l \neq k$, i.e. the challenger has put his entire trust on k . If, in this case, “the playground chosen by the offender” is indeed k , then it is the challenger who should earn points in this duel. More precisely, if $0 < u < 1$, then $P(u, 1) < 0$, and it follows that

$$a > 0. \quad (17)$$

What is then the natural choice for the constant a ? This is now just the matter of normalization. Suppose for the moment that all offenders have chosen playground k . In that case the challenger would receive in total⁶ $-a \cdot \text{card}(R) \cdot P(\bar{x}_k; 1)$ points in the finite case, and $\lim_{n \rightarrow \infty} -a(R_n) \cdot \text{card}(R_n) \cdot P(\bar{x}_k; 1)$ points in the infinite case. It is natural to normalize so that the total is $-P(\bar{x}_k; 1)$ points. Hence we define the constant a to be

$$a = \frac{1}{\text{card}(R)} \quad \text{in the finite case, or} \\ a(R_n) = \frac{1}{\text{card}(R_n)} \quad \text{in the infinite case.} \quad (18)$$

Theorem 1. *If the scoring system satisfies Conditions 1-7 and condition (18), then the resulting system is the Bayesian Truth Serum algorithm, i.e., u^r satisfies (4).*

Proof. Without loss of generality we present the proof for the finite case. In the infinite case we can use exactly the same proof under the limit sign $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{s \in R_n}$.

Using (12) and the Corollary, we obtain

$$u^r = u^r(X, Y) = \sum_{s \in R} T^{s \rightarrow r}(X; Y) - \sum_{s \in R} T^{r \rightarrow s}(X; Y) = \\ = \sum_{s \in R} \sum_{k=1}^m x_k^r \frac{1}{\text{card}(R)} \left(\ln \frac{\bar{x}_k}{y_k^s} \right) - \sum_{s \in R} \sum_{k=1}^m x_k^s \frac{1}{\text{card}(R)} \left(\ln \frac{\bar{x}_k}{y_k^r} \right).$$

The first sum becomes

$$\sum_{s \in R} \sum_{k=1}^m x_k^r \frac{1}{\text{card}(R)} (\ln(\bar{x}_k) - \ln(y_k^s)) = \\ = \sum_{k=1}^m x_k^r \left[\frac{1}{\text{card}(R)} \sum_{s \in R} \ln(\bar{x}_k) - \frac{1}{\text{card}(R)} \sum_{s \in R} \ln(y_k^s) \right].$$

Since the choice of k depends on r (not on s), we obtain

$$\frac{1}{\text{card}(R)} \sum_{s \in R} \ln(\bar{x}_k) = \ln(\bar{x}_k).$$

⁶ In total here means from all the offenders.

On the other hand,

$$\frac{1}{\text{card}(R)} \sum_{s \in R} \ln(y_k^s) = \ln(\widehat{y}_k).$$

It follows that the first sum equals $\sum_{k=1}^m x_k^r \ln\left(\frac{\overline{x}_k}{\widehat{y}_k}\right)$, i.e. equals the information score in (4). For the second sum we obtain

$$\begin{aligned} - \sum_{s \in R} \sum_{k=1}^m x_k^s \frac{1}{\text{card}(R)} \left(\ln \frac{\overline{x}_k}{y_k^s} \right) &= \sum_{s \in R} \sum_{k=1}^m x_k^s \frac{1}{\text{card}(R)} \ln \frac{y_k^r}{\overline{x}_k} = \\ &= \sum_{k=1}^m \ln \frac{y_k^r}{\overline{x}_k} \left(\frac{1}{\text{card}(R)} \sum_{k=1}^m x_k^s \right) = \sum_{k=1}^m \overline{x}_k \ln \frac{y_k^r}{\overline{x}_k}. \end{aligned}$$

This is equal to prediction score in (4).

Q.E.D.

Remark: We would like to emphasize a parallelism between "entropy \leftrightarrow information" vs. "BTS \leftrightarrow information/prediction". This parallelism does not mean that one can be constructed from the other. At this point we are not aware of any approach that axiomatically produces BTS from entropy or vice versa. Perhaps this could be an interesting problem to consider.

First, observe that entropy can be constructed in a similar way, as the one described in this paper. Instead of $(X; Y)$ data, consider only (X) . Instead of playing duels both ways, consider r only as a "challenger" (one can think of it as r "collecting" information data from other players). Hence

$$u^r(X) = - \sum_{s \in R} T^{r \rightarrow s}(X).$$

Suppose that transfers, now only functions of X , satisfy the conditions analogous to the first six conditions in this paper, i.e., we end up with a function $P(x)$. Impose the last condition on P to be the usual additivity condition. Using the same calculation as in the proof of the previous theorem, we obtain that

$$u^r(X) = - \sum_{k=1}^m \overline{x}_k \ln(\overline{x}_k),$$

which is the entropy of \overline{X} . The difference between the input data, i.e. (X, Y) vs. only (X) , is a crucial one. Consider the BTS with the case where Y "does not reveal anything new". More precisely, $y_k^r = \frac{1}{m}$ for every k and r (this, of course, is only for academic purpose). It is then easy to check that, with $x_{k_0}^r = 1$ for a particular k_0 and r , $BTS^r = \text{entropy}(\overline{x}_k) + \ln(x_{k_0})$. Observe that the correction factor $\ln(x_{k_0})$ is precisely the one required to keep the zero sum game property.

Secondly, it is also possible to connect entropy somewhat more directly with the BTS in the following way. From the six

conditions we obtain the form $P(X, Y)$. Assume that we can separate the variables; say $P(x, y) = H(x) - G(y)$. Impose a natural condition that "prediction = actual information" is a draw, i.e., that $P(a, a) = 0$. Obviously then $H = G$. Imposing any entropy-like condition on the second sum (it could be the additivity of G , the proper scoring rule, or even the truth-incentive condition if one wants to work within the Bayesian framework), it can be shown that G is the log function (up to a linear transformation). Consequently, as in the proof of the theorem it follows that $u^r(X) = BTS^r$ (up to a linear transformation).

IV. CONCLUSION

The Bayesian truth serum has been successfully tested on human subjects and in a variety of settings in terms of incentive-compatibility for truth-telling. However, there are situations where telling the truth is not a major issue, but the ranking system is. Moreover, BTS can also be applied in contexts where players are machines (for example measuring information-prediction capability in meteorology, finance, medicine, etc.). In those cases the implementation would shift from truth-telling to ranking systems.

Our ranking is based on a new deterministic mechanism called a "game of duels." There is a large subfamily of those mechanisms in which the ranking of players in the binary case is essentially equivalent to a property we call Gen, which, in the case of infinitely many players, is equivalent to ranking by posterior probabilities. This is similar to the information-cost analysis in which there are many families of functions that will fulfill most properties of standard entropy (so called "sub-exponential" functions can be used instead of entropy; see [21] and the references therein). However, if one wants the additivity property of the uncertainty measure (see [20] for details), then one ends with the standard entropy. Similarly, if one wants additivity for the transfer of points in the game of duels, one ends up with BTS. In future research, it would be of interest to study whether additivity can be replaced by incentive compatibility in a stochastic setting with infinitely many players without additional assumptions that we impose.

REFERENCES

- [1] D. Prelec, "A Bayesian Truth Serum for Subjective Data". Science 306, 462–466, 2004. [\[1\]](#)
- [2] S.R. Miller, B.P. Bailey and A. Kirlik "Exploring the Utility of Bayesian Truth Serum for Assessing Design Knowledge", Human-Computer Interaction, vol.29, no.5-6, pp. 487-515, 2014.

- [3] T. Loughran; R. Paternoster; and K. Thomas, "Incentivizing Responses to Self-report Questions in Perceptual Deterrence Studies: An Investigation of the Validity of Deterrence Theory Using Bayesian Truth Serum". *Journal of Quantitative Criminology*, Vol. 30 Issue 4, p677-707, Dec 2014.
- [4] A. Kukla-Gryz, J. Tyrowicz, M. Krawczyk, K. Siwinski, "We All Do It, but Are We Willing to Admit? Incentivizing Digital Pirates' Confessions", *Applied Economics Letters*, v. 22, iss. 1-3, pp. 184-88, February 2015.
- [5] P. J. Howie., Y. Wang; J. Tsai,. "Predicting new product adoption using Bayesian truth serum", *Journal of Medical Marketing*. Vol. 11 Issue 1, p6-16. 11p, Feb 2011.
- [6] J. Cvitanić, D. Prelec, S. Radas, and H. Šikić, "Mechanism design for an agnostic planner: universal mechanisms, logarithmic equilibrium payoffs and implementation", working paper, October 2015.
- [7] E. Shuford ; A. Albert, and H. E. Massengill, Admissible probability measurement procedures, *Psychometrika*, 31, (2), 125-145, 1966.
- [8] L. J. Savage, "Elicitation of personal probabilities and expectations." *Journal of the American Statistical Association*, 66(336), 783–801, 1971.
- [9] J.M. Bernardo, "Expected Information as Expected Utility". *Ann. Stat.* 7, 686–690, 1978.
- [10] T. Gneiting, and A.E. Raftery, "Strictly Proper Scoring Rules, Prediction, and Estimation". *Journal of the American Statistical Association* 102, 359–378, 2007.
- [11] A.D. Hendrickson, and R.J. Buehler, "Proper scores for probability forecasters". *Ann. Math. Statist.* 42, 1916–1921, 1971.
- [12] A. Banerjee, S. Merugu, I. S. Dhillon, and J. Ghosh. "Clustering with Bregman divergences". *The Journal of Machine Learning Research*, 6:1705–1749, 2005.
- [13] P. Harremoës "Divergence and Sufficiency for Convex Optimization." *Entropy* 19, no. 5: 206, 2017.
- [14] E.Y. Ovcharov (2015) "Proper Scoring Rules and Bregman Divergences". arXiv, arXiv:1502.01178.
- [15] J. Witkowski, and D. C. Parkes, "A Robust Bayesian Truth Serum for Small Populations". In *Proceedings of the 26th AAAI Conference on Artificial Intelligence (AAAI'12)*, 1492–1498, 2012.
- [16] J. Witkowski, "Robust Peer Prediction Mechanisms". Ph.D. Dissertation, Department of Computer Science, Albert-Ludwigs Universität Freiburg. 2014.
- [17] J. Witkowski and D. C. Parkes. "Learning the prior in minimal peer prediction". In *Proceedings of the 3rd Workshop on Social Computing and User Generated Content (SC'13)*, 2013.
- [18] D. Prelec., H.S. Seung, and J. McCoy, "A 'surprisingly popular' solution to the single question crowd wisdom problem". *Nature* 541, 532–535, 2017.
- [19] R. Ray, and D. Prelec. "Creating truth-telling incentives with the Bayesian truth serum." *Journal of Marketing Research* 50, no. 3 , 289-302, 2013.
- [20] R. B. Ash, "Information theory", Dover Publications, 1990.
- [21] H. Šikić and M. V. Wickerhauser, Information cost functions, *Appl. Comput. Harmonic Anal.* 11 147-166. 27, 2001.